

(19) World Intellectual Property  
Organization  
International Bureau



(43) International Publication Date  
31 December 2003 (31.12.2003)

PCT

(10) International Publication Number  
**WO 2004/001623 A2**

(51) International Patent Classification<sup>7</sup>: **G06F 17/20**

[RO/US]; 10970 Palms Blvd. #12, Los Angeles, CA 90034 (US). **KOEHN, Philipp** [DE/US]; Clubhouse Ave #17, Venice, CA 90291 (US).

(21) International Application Number:  
PCT/US2003/009573

(22) International Filing Date: 26 March 2003 (26.03.2003)

(74) Agent: **HARRIS, Scott, C.**; Fish & Richardson P.C., 4350 La Jolla Village Drive, Suite 500, San Diego, CA 92122 (US).

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:  
60/368,070 26 March 2002 (26.03.2002) US  
60/368,447 27 March 2002 (27.03.2002) US

(81) Designated States (*national*): AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NI, NO, NZ, OM, PH, PL, PT, RO, RU, SC, SD, SE, SG, SK, SL, TJ, TM, TN, TR, TT, TZ, UA, UG, US, UZ, VC, VN, YU, ZA, ZM, ZW.

(63) Related by continuation (CON) or continuation-in-part (CIP) to earlier applications:

US 60/368,070 (CIP)  
Filed on 26 March 2002 (26.03.2002)  
US 60/368,447 (CIP)  
Filed on 27 March 2002 (27.03.2002)

(84) Designated States (*regional*): ARIPO patent (GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZM, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HU, IE, IT, LU, MC, NL, PT, RO, SE, SI, SK, TR), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

(71) Applicant (*for all designated States except US*): **UNIVERSITY OF SOUTHERN CALIFORNIA** [US/US]; 3716 South Hope Street, Suite 313, Los Angeles, CA 90007-4344 (US).

**Published:**

— *without international search report and to be republished upon receipt of that report*

(72) Inventors; and

(75) Inventors/Applicants (*for US only*): **MARCU, Daniel** [CA/US]; 2516 Ozone Ct, Hermosa Beach, CA 90254 (US). **KNIGHT, Kevin** [US/US]; 43 6th Street, Hermosa Beach, CA 90245 (US). **MUNTEANU, Dragos, Stefan**

*For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.*

(54) Title: CONSTRUCTING A TRANSLATION LEXICON FROM COMPARABLE, NON-PARALLEL CORPORA

(57) Abstract: A machine translation system may use non-parallel monolingual corpora to generate a translation lexicon. The system may identify identically spelled words in the two corpora, and use them as a seed lexicon. The system may use various clues, e.g., context and frequency, to identify and score other possible translation pairs, using the seed lexicon as a basis. An alternative system may use a small bilingual lexicon in addition to non-parallel corpora to learn translations of unknown words and to generate a parallel corpus.



WO 2004/001623 A2

# CONSTRUCTING A TRANSLATION LEXICON FROM COMPARABLE, NON-PARALLEL CORPORA

## CROSS-REFERENCE TO RELATED APPLICATIONS

[0001] This application claims priority to U.S. Provisional Application Serial No. 60/368,070, filed on March 26, 2002, and U.S. Provisional Application Serial No. 60/368,447, filed on March 27, 2002, the disclosures of which are incorporated by reference.

## ORIGIN OF INVENTION

[0002] The research and development described in this application were supported by DARPA under grant number N66001-00-1-8914. The U.S. Government may have certain rights in the claimed inventions.

## BACKGROUND

[0003] Machine translation (MT) concerns the automatic translation of natural language sentences from a first language (e.g., French) into another language (e.g., English). Systems that perform MT techniques are said to "decode" the source language into the target language.

[0004] Roughly speaking, statistical machine translation (SMT) divides the task of translation into two steps: a word-level translation model and a model for word reordering during the translation process. The statistical models may be trained on parallel corpora. Parallel corpora contain large amounts of text in one language along with their translation in another. Unfortunately, such corpora are available only in limited amounts and cover only in specific genres (Canadian politics, Hong Kong laws, etc). However, monolingual texts exist in higher quantities and in many domains and languages. The availability of monolingual corpora has been enhanced greatly due to the digital revolution and wide-spread use of

the World Wide Web. Methods for processing such resources can therefore greatly benefit the field.

### SUMMARY

**[0005]** In an embodiment, a system may be able to build a translation lexicon from comparable, non-parallel corpora. The system may identify all identically spelled words in the corpora and use these as a seed lexicon for other processes based on clues indicating possible translations.

**[0006]** In another embodiment, a system may align text segments in comparable, non-parallel corpora, matching strings in the corpora, and using the matched strings to build a parallel corpus. The system may build a Bilingual Suffix Tree (BST) and traverse edges of the BST to identify matched strings. The BST may also identify potential translations based on words in the corpora between matched strings.

### BRIEF DESCRIPTION OF THE DRAWINGS

**[0007]** Figure 1 is a block diagram of a system for building a translation lexicon according to an embodiment.

**[0008]** Figure 2 is a flowchart describing a method for building a translation lexicon from non-parallel corpora.

**[0009]** Figure 3 is a table showing results of an experiment utilizing the system of Figure 1.

**[0010]** Figure 4 is a block diagram of a system for building a translation lexicon according to another embodiment.

**[0011]** Figure 5 is a suffix tree.

**[0012]** Figure 6 is a Generalized Suffix Tree (GST).

**[0013]** Figure 7 is a Bilingual Suffix Tree (BST).

**[0014]** Figure 8 is a portion of a BST showing example alignments.

**[0015]** Figure 9 are portions of a BST describing left and right alignments.

[0016] Figure 10 is pseudocode describing an algorithm for learning translations of unknown words.

#### DETAILED DESCRIPTION

[0017] Figure 1 shows a system 100 for building a translation lexicon 105 according to an embodiment. The system may use non-parallel monolingual corpora 110, 115 in two languages to automatically generate one-to-one mapping of words in the two languages.

[0018] The two monolingual corpora should be in a fairly comparable domain. For example, in an implementation, an English-German translation lexicon was generated from a 1990-1992 Wall Street Journal corpus on the English side and a 1995-1996 German news wire (DPA) on the German side. Both corpora are news sources in the general sense. However, they span different time periods and have a different orientation: the World Street Journal covers mostly business news, the German news wire mostly German politics.

[0019] The system 100 may use clues to find translations of words in the monolingual corpora. The first clue considered may be the existence of identical words in the two corpora. Due to cultural exchange, a large number of words that originate in one language may be adopted by others. Recently, this phenomenon can be seen with words such as "Internet" or "Aids". These terms may be adopted verbatim or changed by well-established rules. For instance, "immigration" (German and English) has the Portuguese translation "imigração", as many words ending in -tion have translations with the same spelling except for the ending changed to -ção.

[0020] Figure 2 shows a flowchart describing a method 200 for building a translation lexicon from non-parallel corpora. A word comparator 120 may be used to collect pairs of identical words (block 205). In the English-German implementation described above, 977 identical words were

found. When checked against a benchmark lexicon, the mappings were found to be 88% correct.

**[0021]** The correctness of word mappings acquired in this fashion may depend highly on word length. While identical three-letter words were only translations of each other 60% of the time, this was true for 98% of ten-letter words. Clearly, for shorter words, the accidental existence of an identically spelled word in the other language word is much higher. Accordingly, the word comparator 120 may restrict the word length to be able to increase the accuracy of the collected word pairs. For instance, by relying only on words at least of length six, 622 word pairs were collected with 96% accuracy.

**[0022]** The identified identically spelled word pairs may be used as a seed lexicon 130 (block 210). A lexicon builder 125 may expand the seed lexicon into the larger translation lexicon 105 by applying rules based on clues which indicate probable translations. The lexicon builder 125 may use seed lexicon to bootstrap these methods, using the word pairs in the seed lexicon as correct translations.

**[0023]** As already mentioned, there are some well-established transformation rules for the adoption of words from a foreign language. For German to English, this includes replacing the letters k and z by c and changing the ending -tät to -ty. Both these rules can be observed in the word pair Elektrizität and electricity. The lexicon builder 125 may utilize these rules to expand the seed lexicon. In the English-German implementation, 363 additional word pairs were collected, with an accuracy of 91%.

**[0024]** The lexicon builder 125 extracts potential translation word pairs based on one or more clues. These clues may include similar spelling, similar context, preserving word similarity, and word frequency.

**[0025]** When words are adopted into another language, their spelling might change slightly in a manner that cannot be simply generalized in a rule, e.g., "website" and "Webseite". This is even more the case for words that can be traced back to common language roots, such as "friend" and "Freund", or "president" and "Präsident". Still, these words, often called "cognates", maintain a very similar spelling. This can be defined as differing in very few letters. This measurement can be formalized as the number of letters common in sequence between the two words, divided by the length of the longer word.

**[0026]** The example word pair "friend" and "freund" shares 5 letters (fr-e-nd), and both words have length 6, hence their spelling similarity is  $5/6$ , or 0.83. This measurement may be referred to as the "longest common subsequence ratio." The lexicon builder 125 may measure the spelling similarity between every German and English word, and sort possible word pairs accordingly. This may be done in a greedy fashion, i.e., once a word is assigned to a word pair, the lexicon builder 125 does not look for another match.

**[0027]** Another clue is context. If the monolingual corpora are somewhat comparable, it can be assumed that a word that occurs in a certain context should have a translation that occurs in a similar context. The context may be defined by the frequencies of context words in surrounding positions. This context has to be translated into the other language, and the lexicon builder 125 can search the word with the most similar context.

**[0028]** The lexicon builder 125 may collect counts over words occurring in an n-word window, e.g., four words ( $n=4$ ), around the target word. For each occurrence of a target word, the counts may be collected over how often certain context words occur in the two positions directly ahead of the target word and the two following positions. The counts may be

collected separately for each position and then entered into a context vector with a dimension for each context word in each position. Finally, the raw counts are normalized. Vector comparison is done by adding all absolute differences of all components.

**[0029]** Alternatively, the lexicon builder 125 may count how often another word occurs in the same sentence as the target word. The counts may then be normalized by a using the tf/idf method, which is often used in information retrieval.

**[0030]** The seed lexicon may be used to construct context vectors that contain information about how a new unmapped word co-occurs with the seed words. This vector can be translated into the other language, since we already know the translations of the seed words are already known. The lexicon builder 125 can search for the best matching context vector in the target language, and decide upon the corresponding word to construct a word mapping. The lexicon builder 125 may compute all possible word, or context vector, matches. The best word matches may be collected in a greedy fashion.

**[0031]** Another clue is based on the assumption that pairs of words that are similar in one language should have translations that are similar in the other language. For instance, Wednesday is similar to Thursday as Mittwoch is similar to Donnerstag. Two words may be defined as similar if they occur in a similar context, which is the case for Wednesday and Thursday.

**[0032]** In one approach, the context vector for each word in the lexicon may consist of co-occurrence counts in respect to a number of peripheral tokens (basically, the most frequent words). These counts may be collected for each position in an n-word window around the word in focus.

**[0033]** Instead of comparing the co-occurrence counts directly, the Spearman rank order correlation may be applied. For each position, the tokens are compared in frequency and

the frequency count is replaced by the frequency rank, e.g., the most frequent token count is replaced with 1 and the least frequent by  $n$ . The similarity of the two context vectors  $a = (a_i)$  and  $b = (b_i)$  is then defined by:

$$[0034] \quad R(a,b) = 1 - \frac{6\sum(a_i - b_i)}{4n(n^2 - 1)}$$

[0035] The result is a matrix with similarity scores between all German words, and a second matrix with similarity scores between all English words. For a new word, the lexicon builder 125 may look up its similarity scores to seed words, thus creating a similarity vector. Such a vector can be translated into the other language. The translated vector can be compared to other vectors in the second language.

[0036] The lexicon builder 125 may perform a greedy search for the best matching similarity vectors and add the corresponding words to the lexicon.

[0037] Another clue is based on the assumption that in comparable corpora, the same concepts should occur with similar frequencies. Frequency may be defined as a ratio of the word frequencies normalized by the corpus sizes.

[0038] Each of the clues provides a matching score between two words (block 220), e.g., a German word and an English word. The likelihood of these two words being actual translations of each other may correlate to these scores. The lexicon builder 125 may employ a greedy search to determine the best set of lexicon entries based on these scores (block 225). First, the lexicon builder 125 searches for the highest score for any word pair. This is added to the lexicon (block 230), and word pairs that include either the German and English word are dropped from further search. This may be performed iteratively until all words are used up.

[0039] The lexicon builder 125 may combine different clues by adding up the matching scores. The scores can be weighted. For example, when using the spelling clue in combination with



others, it may be useful to define a cutoff. If two words agree in 30% of their letters, this is generally as bad as if they do not agree in any, i.e., the agreements are purely coincidental.

**[0040]** Figure 3 shows results of the English-German implementation. "Entries" indicate the number of correct lexicon entries that were added to a seed lexicon of 1337 identically spelled words, and "Corpus" indicates how well the resulting translation lexicon performs compared to the actual word-level translations in a parallel corpus.

**[0041]** The English-German implementation was restricted to nouns. Verbs, adjectives, adverbs and other part of speech may be handled in a similar way. They might also provide useful context information that is beneficial to building a noun lexicon. These methods may be also useful given a different starting point. For example, when building machine translation systems, some small parallel text should be available. From these, some high-quality lexical entries can be learned, but there will always be many words that are missing. These may be learned using the described methods.

**[0042]** Figure 4 shows a system 400 for building a translation lexicon according to another embodiment. The system 400 may also be used to build parallel corpora from comparable corpora. Given an initial bilingual lexicon 405 and two texts 410, 415 in each of the languages, the system 400 may identify parts of the texts which can be aligned (i.e., are mutual translations of each other according to the lexicon). The parts can be arbitrarily long, i.e., the system 400 may align sequences of a few words rather than or in addition to whole sentences or whole phrases. Based on these alignments, the system 400 may generate a parallel corpus 420 and identify translations 425 of words from the source language which are not in the lexicon.

[0043] For example, consider the following two sentences where the only unknown French word is "raison":

"Ce est pour cette raison que le initiative de le ministre..."; and

"It is for this reason that the party has proposed..."

Since "Ce est pour cette" can be aligned with "It is for this" and "que le" with "that the", it is a reasonable assumption that "raison" can be translated by "reason". The system 400 may search the corpora for cases similar to this example.

[0044] The system 400 may use a suffix tree data structure in order to identify the alignments. The suffix tree of a string uniquely encodes all the suffixes of that string (and thus, implicitly, all its substrings too). The system 400 may first build such a tree of the target language corpus, and then add to each substrings all the substrings from the source language corpus that align to it. The next step is to identify unknown target language words that are surrounded by aligned substrings. The source language word that corresponds to the "well-aligned" unknown is considered to be a possible translation.

[0045] A suffix tree stores in linear space all suffixes of a given string. Such succinct encoding exposes the internal structure of the string, providing efficient (usually linear-time) solutions for many complex string problems, such as exact and approximate string matching, finding the longest common substring of multiple strings, and string compression. Formally, a suffix tree for a string S of length N has the following properties: each edge of the tree is labeled by a nonempty substring of S; each internal node other than the root has at least two children; no two edges out of a node can have edge-labels beginning with the same character/word; and (the key feature of the tree) there is a one-to-one correspondence between all suffixes of S and paths in the tree from the root to the leaves.

**[0046]** Figure 5 shows the suffix tree 500 of string xyzyxzy. Note that if a suffix of a string is also a prefix of another suffix (as would be the case for suffix zy of string xyzyxzy), a proper suffix tree cannot be built for the string. The problem is that the path corresponding to that suffix would not end at a leaf, so the tree cannot have the last property in the list above. To avoid this, the system 400 appends an end-of-string marker "\$" that appears nowhere else in the string. For clarity, the drawings only show the \$ marker when necessary.

**[0047]** Each monolingual corpus given as input to the system 400 may be divided into a set of sentences. The system 400 may use a variant of suffix trees that works with sets of strings, namely Generalized Suffix Trees (GST). In a GST of a set of strings, each path from the root to a leaf represents a suffix in one or more strings from the set. A conceptually easy way to build such a tree is to start by building a regular suffix tree for the first sentence in the corpus, and then for each of the other sentences to take their suffixes one by one and add them to the tree (if they are not already in it). Figure 6 shows the GST 600 for a corpus of two sentences. The numbers at the leaves 605 of the tree show which sentences contain the suffix that ends there.

**[0048]** Building the suffix tree of a string takes time and space linear in the length of the string. Building a GST for a set of strings takes time and space linear in the sum of the lengths of all strings in the set.

**[0049]** A Bilingual Suffix Tree (BST) is the result of matching a source language GST against a target language GST. Two strings (i.e., sequences of words) match if the corresponding words are translations of each other according to a bilingual lexicon. In order to perform the matching operation, all paths that correspond to an exhaustive traversal of one of the trees (the source tree) are traversed

in the other (the target tree), until a mismatch occurs. In the process, the target tree is augmented with information about the alignments between its paths and those of the source, thus becoming a bilingual suffix tree. Figure 7 shows two corpora 705, 710, a bilingual lexicon 715, and the corresponding BST 720. Edges drawn with dotted lines mark ends of alignment paths through the tree. Their labels are (unaligned) continuations of the source language substrings from the respective paths.

**[0050]** Since there is a one-to-one correspondence between the substrings in the text and the paths in the suffix trees, the operation described above will yield all pairs of substrings in the two corpora given as input and discover all partial monotone alignments defined by the lexicon.

**[0051]** If the lexicon is probabilistic, each matching between two words will be weighted by the corresponding translation probability. The paths in the resulting bilingual tree will also have weights associated with them, defined as the product of the matching probabilities of the words along the path.

**[0052]** BSTs are constructed to encode alignment information, therefore the extraction of parallel phrases amounts to a simple depth-first traversal of the tree. Figure 8 shows some alignments we can extract from the BST in Figure 7, a portion of which is shown in Figure 8.

**[0053]** As can be seen in Figure 4, there are three types of edge labels in a BST: only target language sequences (e.g., xzy), pairs of target and source language sequences (y:b followed by z:c) and only source language words (b or c). For alignment extraction we are interested in edges of the third type, because they mark ends of alignments. Let e be an edge labeled only with a source language word, originating from node n. A path from the root to n will only traverse edges labeled with word pairs, defining two aligned sequences. The

fact that  $n$  has outgoing edge  $e$  indicates there is a mismatch on the subsequent words of those two sequences. Thus, in order to extract all aligned substrings, the system 400 traverses the BST on edges labeled with word pairs, and extract all paths that end either at the leaves or at nodes that have outgoing edges labeled only with source language words.

**[0054]** The heuristic by which the system 400 discovers new word translations is shown graphically in Figure 9. Figure 9(i) shows a branch 905 of the BST corresponding to the comparable corpus in the same figure. The path defined by the bold edges shows that sequences  $xyz$  and  $abc$  are aligned, and diverge (i.e., have a mismatch) at characters  $y$  and  $d$  respectively. This may be taken as a weak indication that  $d$  and  $y$  are translations of each other. This indication would become stronger if, for example, the sequences following  $d$  and  $y$  in the two corpora would also be aligned. One way to verify this is to reverse both strings, build a BST for the reversed corpora (a reverse BST), and look for a common path that diverges at the same  $d$  and  $y$ . Figure 9(ii) shows the reverse BST 910, and in bold, the path we are interested in. When  $d$  and  $y$  are surrounded by aligned sequences, we hypothesize that they are translations of each other.

**[0055]** For a pair of words from the two corpora, we use the terms "right alignment" and "left alignment" to refer to the aligned sequences that precede and respectively succeed the two words in each corpus. The left and right alignments and the two words delimited by them make up a context alignment. For example, the left alignment  $xyz-abc$ , the right alignment  $xzy-acb$  and the words  $y$  and  $d$  in Figure 9(iii) make up a context alignment 915.

**[0056]** Given a comparable corpus, this procedure will yield many context alignments which correspond to incorrect

translations, such as that between the words "canadien" and "previous":

**tout canadien sérieux**

**any previous serious**

**[0057]** In order to filter out such cases, the system 400 uses two simple heuristics: length and word content. Thus, for a context alignment to be valid, the left and right context together must contain at least three words, one of which must be an open-class word, e.g., a noun, verb, adjective, or adverb, classes which can have new words added to them. The translation candidate must also be an open-class word. The algorithm 1000 for learning translations of unknown words is summarized in Figure 10. An advantage of the algorithm over previous approaches is that we do not provide as input to the algorithm a list of unknown words. Instead, the system automatically learns from the corpus both the unknown words and their translation, upon discovery of appropriate context alignments.

**[0058]** The system 400 was tested on an English-French comparable corpus, of approximately 1.3 million words -50.000 sentences for each language. It was obtained by taking two non-parallel, nonaligned segments from the Hansard corpus. The Hansard Corpus includes parallel texts in English and Canadian French, drawn from official records of the proceedings of the Canadian Parliament. A small bilingual lexicon of 6,900 entries was built using 5,000 sentences pairs (150,000 words for each language). The parallel corpus was taken from the Proceedings of the European Parliament (EuroParl) Note that the parallel corpus belongs to a different domain than the comparable corpus. Also the parallel corpus is extremely small. For low density languages, such a corpus can be built manually.

**[0059]** When given as input the comparable corpora described above and the bilingual lexicon of 6,900 entries, the

algorithm 1000 found 33,926 parallel sequences, with length between three and seven words. Out of 100 randomly selected sequences, 95% were judged to be correct.

[0060] The system also found translations for thirty unknown French words. Of these, nine were correct, which means a precision of 30%.

[0061] For each of the two corpora, building the monolingual GST took only 1.5 minutes. The matching operation that yields the BST is the most time-consuming: it lasted 38 hours for the forward BST and 60 hours for the reverse BST. The extractions of all parallel phrases and of the translations took about 2 hours each. The experiments were run on a Linux® system 400 with an Intel® Pentium® 3 processor of 866Mhz.

[0062] A number of embodiments have been described. Nevertheless, it will be understood that various modifications may be made without departing from the spirit and scope of the invention. For example, blocks in the flowcharts may be skipped or performed out of order and still produce desirable results. Also, the heuristics described in paragraphs 0018-0040 may be combined with the alignment method described in paragraphs 0041-0060. Accordingly, other embodiments are within the scope of the following claims.

**CLAIMS**

1. A method for building a translation lexicon from non-parallel corpora, the method comprising:

identifying identically spelled words in a first corpus and a second corpus, the first corpus including words in a first language and the second corpus including words in a second language;

generating a seed lexicon including identically spelled words; and

expanding the seed lexicon by identifying possible translations of words in the first and second corpora using one or more clues.

2. The method of claim 1, wherein said expanding comprises using the identically spelled words in the seed lexicon as accurate translations.

3. The method of claim 1, further comprising:  
identifying substantially identical words in the first and second corpora; and

adding said substantially identical words to the seed lexicon.

4. The method of claim 3, wherein said identifying substantially identical words comprises applying transformation rules to words in the first corpora to form transformed words; and

comparing said transformed words to words in the second corpora.

5. The method of claim 1, wherein said one or more clues includes similar spelling.



6. The method of claim 1, wherein said identifying comprises identifying cognates.

7. The method of claim 1, wherein said identifying comprises identifying word pairs having a minimum longest common subsequence ratio.

8. The method of claim 1, wherein said one or more clues includes similar context.

9. The method of claim 1, wherein said identifying comprises:

identifying a plurality of context words; and  
identifying a frequency of context words in an n-word window around a target word.

10. The method of claim 9, further comprising generating a context vector.

11. The method of claim 1, wherein said identifying comprises identifying frequencies of occurrence of word in the first and second first corpora.

12. The method of claim 1, further comprising:  
generating matching scores for each of a plurality of clues.

13. The method of claim 12, further comprising adding the matching scores.

14. The method of claim 13, further comprising weighting the matching scores.

15. A method for generating parallel corpora from non-parallel corpora, the method comprising:

aligning text segments in two non-parallel corpora, the corpora including a source language corpus and a target language corpus;

matching strings in the two non-parallel corpora; and  
generating a parallel corpus including the matched strings as translation pairs.

16. The method of claim 15, wherein said matching comprises using a bilingual lexicon comprising translation pairs including corresponding source language words and target language words.

17. The method of claim 15, wherein said aligning comprises:

generating a Bilingual Suffix Tree.

18. The method of claim 17, further comprising:

traversing the Bilingual Suffix Tree on edges labeled with word pairs; and

extracting paths that end at one of a leaf and a node having outgoing edges labeled only with source language words.

19. The method of claim 17, further comprising:

generating a Generalized Suffix Tree from the source language corpus;

generating a Generalized Suffix Tree from the target language corpus; and

matching strings in said Generalized Suffix Trees.

20. The method of claim 15, further comprising:

identifying words in the two corpora surrounded by matching strings, one of the words being unknown.

21. The method of claim 20, further comprising:  
identifying said words as a translation pair.

22. The method of claim 20, further comprising:  
generating a Bilingual Suffix Tree from the two corpora;  
generating a reverse Bilingual Suffix Tree; and  
identifying words in the two corpora surrounded by  
aligned sequences.

23. An apparatus comprising:  
a word comparator operative to identify identically  
spelled words in a first corpus and a second corpus and build  
a seed lexicon including said identically spelled words, the  
first corpus including words in a first language and the  
second corpus including words in a second language; and  
a lexicon builder operative to expand the seed lexicon by  
identifying possible translations of words in the first and  
second corpora using one or more clues.

24. The apparatus of claim 23, wherein use the  
identically spelled words in the seed lexicon as accurate  
translations.

25. An apparatus for generating parallel corpora from  
non-parallel corpora, the apparatus comprising:

an alignment module operative to align text segments in  
two non-parallel corpora, the corpora including a source  
language corpus and a target language corpus; and

a matching module operative to match strings in the two  
non-parallel corpora generate a parallel corpus including the  
matched strings as translation pairs.

26. The apparatus of claim 25, wherein the aligning module is operative to build a Bilingual Suffix Tree from a text segment from one of said two non-parallel corpora.

27. An article comprising a machine-readable medium including machine-executable instructions, the instructions operative to cause a machine to:

identify identically spelled words in a first corpus and a second corpus, the first corpus including words in a first language and the second corpus including words in a second language;

generate a seed lexicon including identically spelled words; and

expand the seed lexicon by identifying possible translations of words in the first and second corpora using one or more clues.

28. The article of claim 27, wherein the instructions operative to cause the machine to expand comprise instructions operative to cause the machine to use the identically spelled words in the seed lexicon as accurate translations.

29. An article comprising a machine-readable medium including machine-executable instructions, the instructions operative to cause a machine to:

align text segments in two non-parallel corpora, the corpora including a source language corpus and a target language corpus;

match strings in the two non-parallel corpora; and

generate a parallel corpus including the matched strings as translation pairs.

30. The article of claim 29, wherein the instructions operative to cause the machine to match comprise instructions

operative to cause the machine to use a bilingual lexicon comprising translation pairs including corresponding source language words and target language words.

1/9

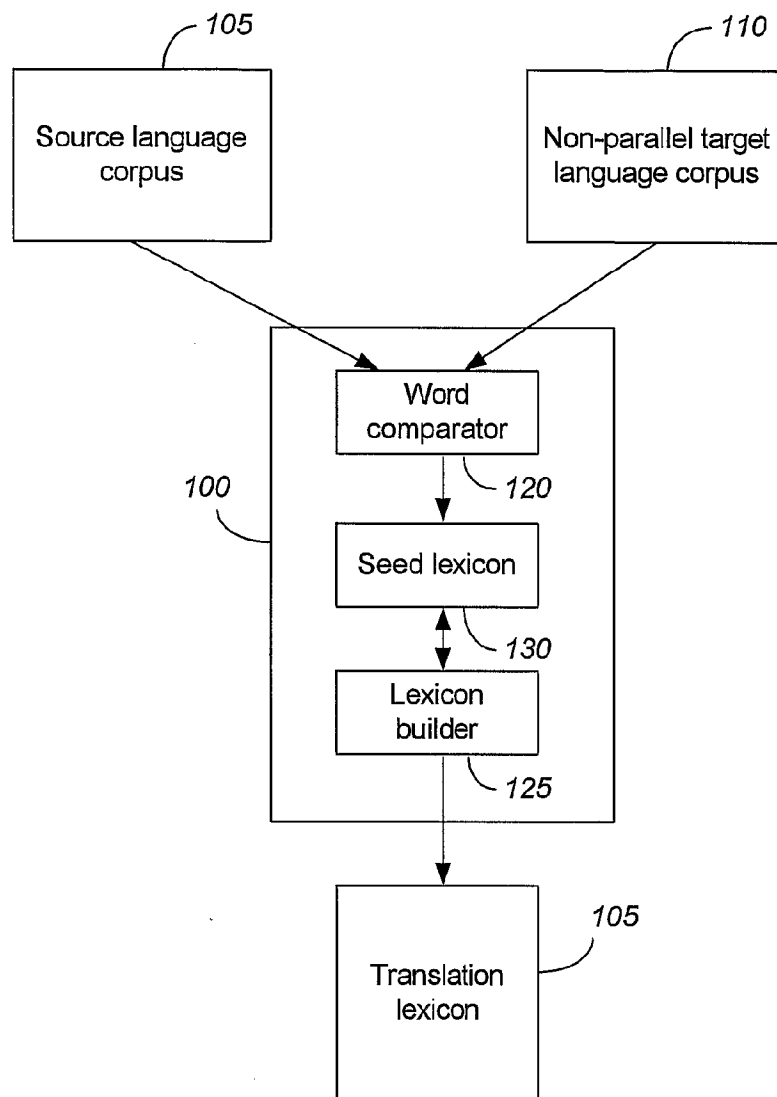


FIG. 1

200 →

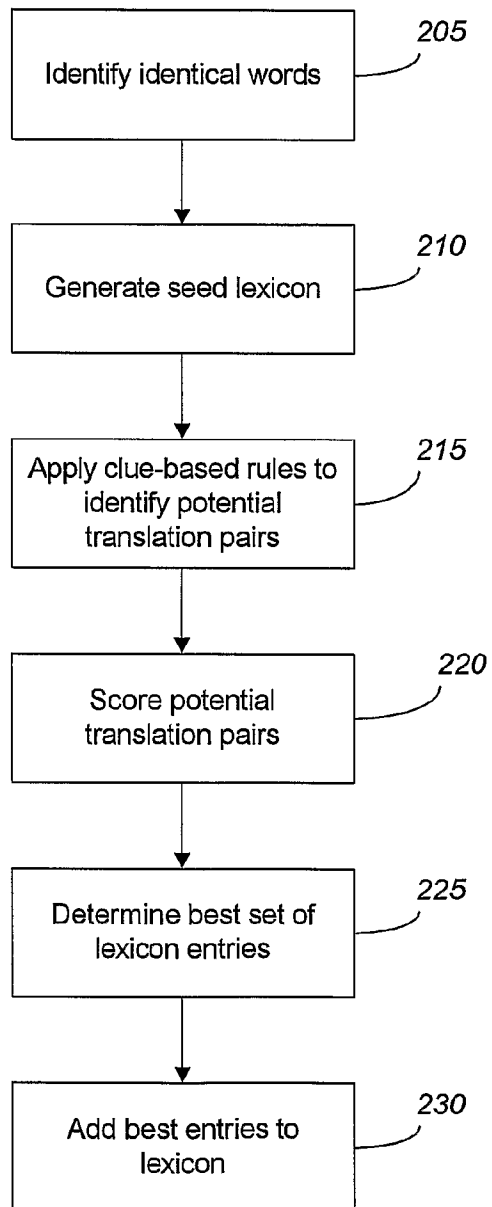


FIG. 2

Clues	Entries	Corpus
Spelling	140	25.4%
Context	107	31.9%
Preserving Similarity	2	15.8%
Frequency	2	17.6%
Spelling+Context	185	38.6%
Spelling+Frequency	151	27.4%
Spelling+Context+Similarity	186	39.0%
All clues	186	39.0%

FIG. 3



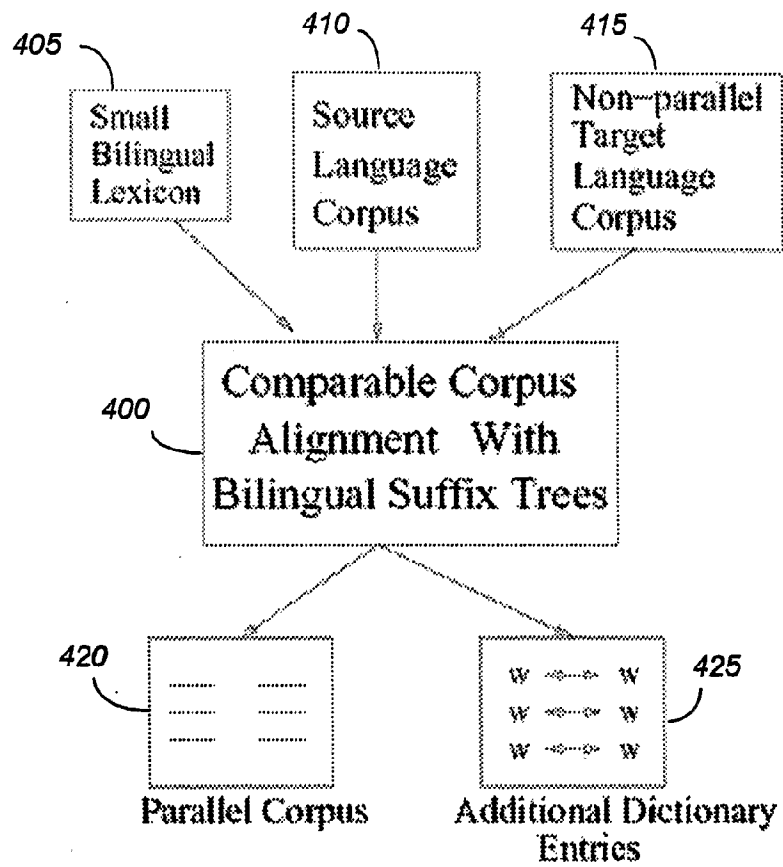


FIG. 4

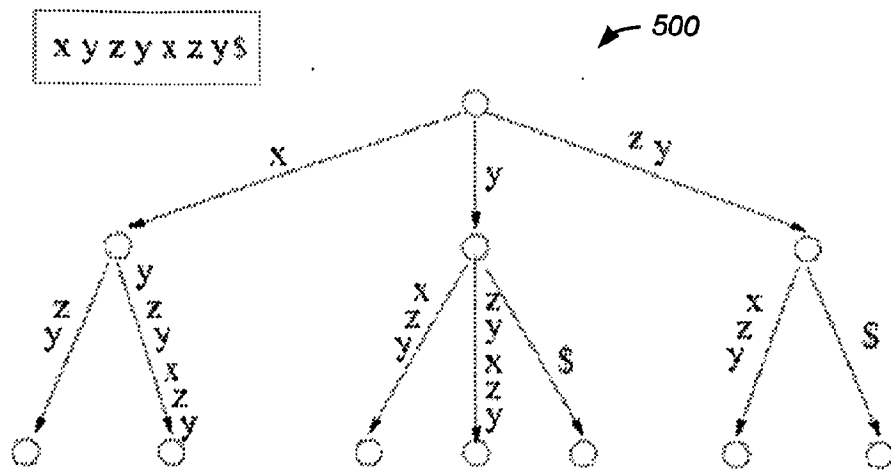


FIG. 5

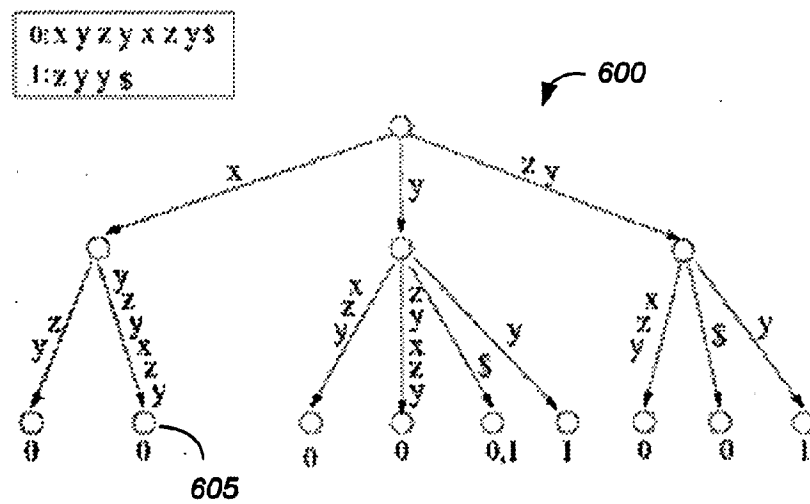


FIG. 6

6/9

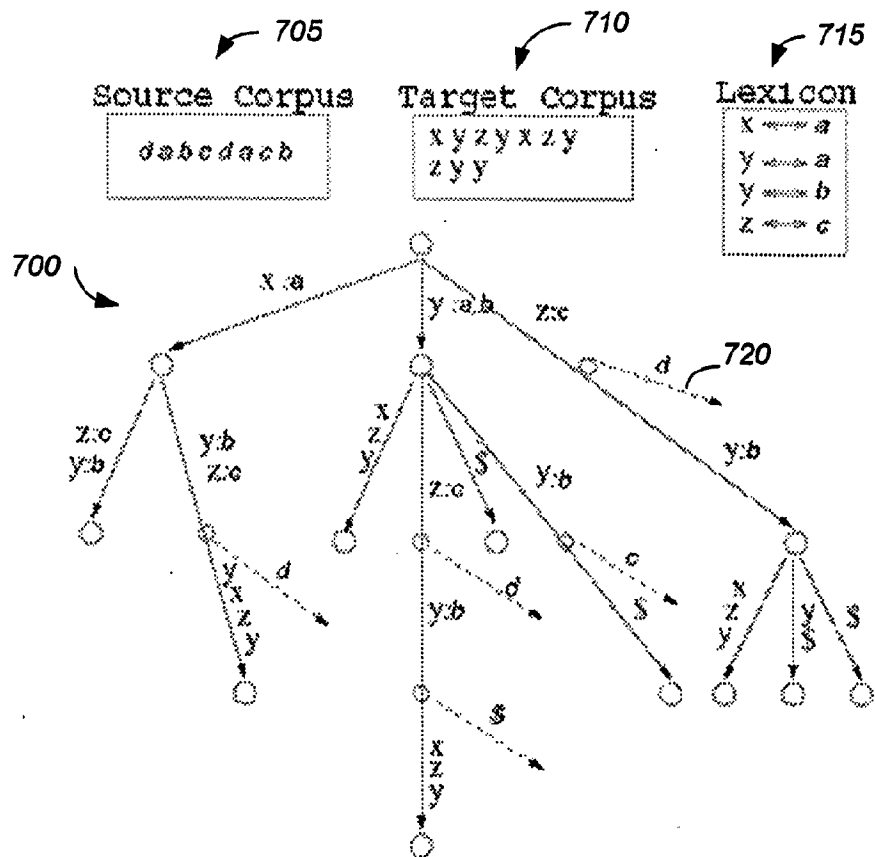


FIG. 7

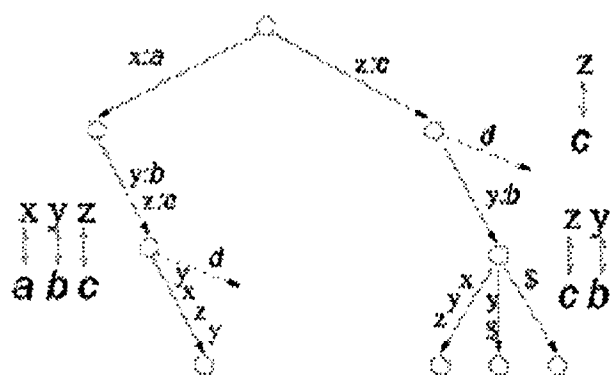


FIG. 8

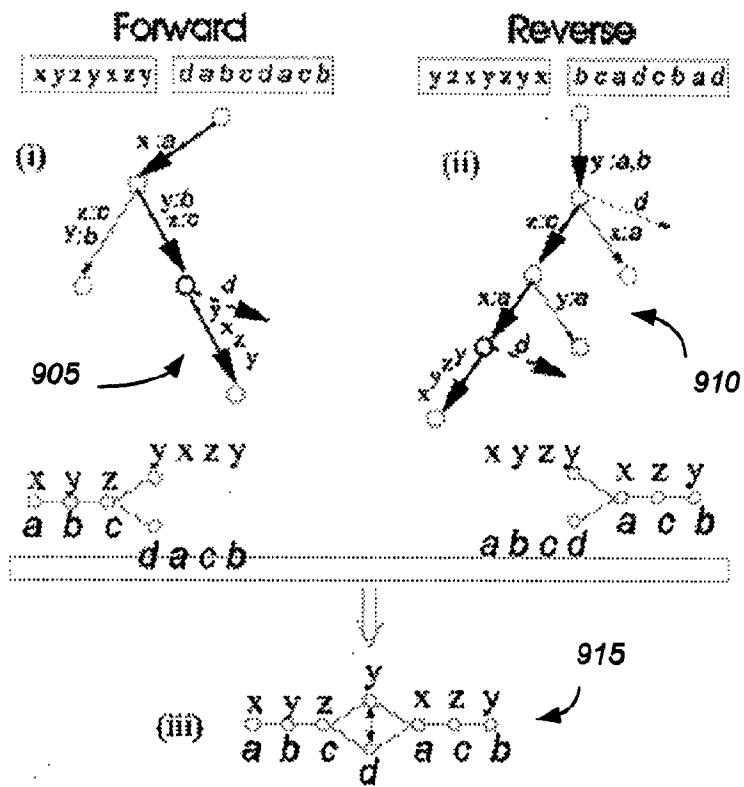



FIG. 9

1000 

1. Build the forward and backward BSTs.
  2. Traverse each BST and extract left and right alignments for every node that represents a divergence
- For each word pair from the divergence set:
- a. create context alignments out of appropriate left and right alignments
  - b. filter out invalid context alignments
  - c. extract valid translation candidates from the context alignments

FIG. 10